

H.W. 2.1.a)

$f$  is Lipschitz continuous iff  $\|w\|_2 \leq D \Rightarrow \|\nabla f(w)\|_2 \leq B$

We have  $f(w) = \frac{1}{N} \sum_{i \in [N]} f_i(w) + \lambda \|w\|_2^2$  where  $f_i(w) = \log(1 + \exp\{-y_i w^T x_i\})$ .

Since  $\nabla f_i(w) = \frac{\exp\{-y_i w^T x_i\}}{1 + \exp\{-y_i w^T x_i\}} (-y_i x_i) = \frac{-y_i x_i}{1 + \exp\{y_i w^T x_i\}}$ , we have:

$$\begin{aligned} \|\nabla f(w)\|_2 &= \left\| \frac{1}{N} \sum_{i \in [N]} \nabla f_i(w) + 2\lambda w \right\|_2 = \left\| \frac{1}{N} \sum_{i \in [N]} \underbrace{\frac{-y_i x_i}{1 + \exp\{y_i w^T x_i\}}}_{\geq 0} + 2\lambda w \right\|_2 \\ &\leq \frac{1}{N} \left\| \sum_{i \in [N]} y_i x_i \right\|_2 + 2\lambda \|w\|_2 \end{aligned}$$

Hence, if  $\|w\|_2 \leq D$ , then  $\|\nabla f(w)\|_2 \leq \frac{1}{N} \left\| \sum_{i \in [N]} y_i x_i \right\|_2 + 2\lambda D$

This shows that  $f$  is Lipschitz continuous with constant

$$B = \frac{1}{N} \left\| \sum_{i \in [N]} y_i x_i \right\|_2 + 2\lambda D.$$

H.W. 2.1 b)

We have from Lecture 3 that:

$f_i(w) = \log(1 + \exp(-y_i w^T x_i))$  is  $L$ -smooth iff  $f_i(w_2) \leq f_i(w_1) + \nabla f_i(w_1)^T (w_2 - w_1) + \frac{L}{2} \|w_2 - w_1\|_2^2$

From lecture 2 & assignment 1.2a), we have that this holds if:  $\|\nabla f_i(w_2) - \nabla f_i(w_1)\|_2 \leq L \|w_2 - w_1\|_2$

$$\text{We have: } \nabla f_i(w) = \frac{\exp\{-y_i w^T x_i\}}{1 + \exp\{-y_i w^T x_i\}} (-y_i x_i) = \frac{(-y_i x_i)}{1 + \exp\{y_i w^T x_i\}}$$

$$\Rightarrow \|\nabla f_i(w_2) - \nabla f_i(w_1)\|_2 = \left\| \frac{(-y_i x_i)}{1 + \exp\{y_i w_2^T x_i\}} - \frac{(-y_i x_i)}{1 + \exp\{y_i w_1^T x_i\}} \right\|_2 \quad (2.1.b.i)$$

$$= \|y_i x_i\|_2 \left\| \frac{1}{1 + \exp\{y_i w_2^T x_i\}} - \frac{1}{1 + \exp\{y_i w_1^T x_i\}} \right\|_2$$

Since  $0 \leq \frac{1}{1 + \exp\{a\}} \leq 1$  for  $a \in \mathbb{R}$ , we have  $\left\| \frac{1}{1 + \exp\{b\}} - \frac{1}{1 + \exp\{a\}} \right\|_2 \leq \|1 - 0\|_2 \leq 1$  for  $a, b \in \mathbb{R}$

Hence, (2.1.b.i) gives:  $\|\nabla f_i(w_2) - \nabla f_i(w_1)\|_2 \leq \|y_i x_i\|_2 \quad (2.1.b.ii)$

Hence,  $f_i(w)$  is  $L$ -smooth with  $L = \frac{\|y_i x_i\|_2}{\|w_2 - w_1\|_2}$

$$\text{Since } f(w) = \frac{1}{N} \sum_{i \in [N]} f_i(w) + \lambda \|w\|_2^2 \Rightarrow \nabla f(w) = \frac{1}{N} \sum_{i \in [N]} \nabla f_i(w) + 2\lambda w$$

$$\text{We have: } \|\nabla f(w_2) - \nabla f(w_1)\|_2 = \left\| \left( \frac{1}{N} \sum_{i \in [N]} \nabla f_i(w_2) + \lambda w_2 \right) - \left( \frac{1}{N} \sum_{i \in [N]} \nabla f_i(w_1) + \lambda w_1 \right) \right\|_2 =$$

$$\leq \frac{1}{N} \sum_{i \in [N]} \|\nabla f_i(w_2) - \nabla f_i(w_1)\|_2 + |\lambda| \|w_2 - w_1\|_2 \stackrel{(2.1.b.ii)}{\leq} \frac{1}{N} \sum_{i \in [N]} \|y_i x_i\|_2 + |\lambda| \|w_2 - w_1\|_2$$

$$\text{Hence, } f(w) \text{ is } L\text{-smooth with } L = \frac{\frac{1}{N} \sum_{i \in [N]} \|y_i x_i\|_2}{\|w_2 - w_1\|_2} + |\lambda|$$

$\nwarrow$  should be positive

# H.W. 2.1 c)

$f$  is strongly convex with constant  $\mu > 0$  iff

$$f(w_2) \geq f(w_1) + \nabla f(w_1)^T (w_2 - w_1) + \frac{\mu}{2} \|w_2 - w_1\|_2^2, \text{ which is equivalent to:}$$

$$g(w) = f(w) - \frac{\mu}{2} \|w\|^2 \text{ is convex because of the 1st order condition for convexity}$$

Following the monotone gradient condition for convexity, we have therefore that

$$f(w) \text{ is strongly convex iff } (\nabla f(w_2) - \nabla f(w_1) - \mu w_2)^T (w_2 - w_1) \geq 0,$$

$$\text{Which is equivalent to } (\nabla f(w_2) - \nabla f(w_1))^T (w_2 - w_1) \geq \mu \|w_2 - w_1\|_2^2$$

$$\text{We have: } \nabla f_i(w) = \frac{\exp\{-y_i w^T x_i\}}{1 + \exp\{-y_i w^T x_i\}} (-y_i x_i) = \frac{-y_i x_i}{1 + \exp\{y_i w^T x_i\}}$$

$$\text{and } \nabla f(w) = \frac{1}{N} \sum_{i \in [N]} \nabla f_i(w) + 2\lambda w$$

$$\begin{aligned} (\nabla f(w_2) - \nabla f(w_1))^T (w_2 - w_1) &= \left[ \left( \frac{1}{N} \sum_{i \in [N]} \nabla f_i(w_2) + 2\lambda w_2 \right) - \left( \frac{1}{N} \sum_{i \in [N]} \nabla f_i(w_1) + 2\lambda w_1 \right) \right]^T (w_2 - w_1) = \\ &= \left( \frac{1}{N} \sum_{i \in [N]} \nabla f_i(w_2) - \nabla f_i(w_1) \right)^T (w_2 - w_1) + 2\lambda (w_2 - w_1)^T (w_2 - w_1) = \\ &= \left( \frac{1}{N} \sum_{i \in [N]} \frac{-y_i x_i}{1 + \exp\{y_i w_2^T x_i\}} - \frac{-y_i x_i}{1 + \exp\{y_i w_1^T x_i\}} \right)^T (w_2 - w_1) + 2\lambda \|w_2 - w_1\|_2^2 = \\ &= \left( \frac{1}{N} \sum_{i \in [N]} \frac{y_i x_i}{1 + \exp\{y_i w_1^T x_i\}} - \frac{y_i x_i}{1 + \exp\{y_i w_2^T x_i\}} \right)^T (w_2 - w_1) + 2\lambda \|w_2 - w_1\|_2^2 = \\ &= \left( \frac{1}{N} \sum_{i \in [N]} \frac{y_i x_i (\exp\{y_i w_2^T x_i\} - \exp\{y_i w_1^T x_i\})}{(1 + \exp\{y_i w_1^T x_i\})(1 + \exp\{y_i w_2^T x_i\})} \right)^T (w_2 - w_1) + 2\lambda \|w_2 - w_1\|_2^2 \geq \\ &\geq 0 + 2\lambda \|w_2 - w_1\|_2^2 = 2\lambda \|w_2 - w_1\|_2^2 \end{aligned}$$

Hence,  $f(w)$  is strongly convex with  $\mu = 2\lambda$ .

H.W. 2.2.

Since  $\text{Var}[X] = \mathbb{E}[X^2] - \mathbb{E}[X]^2 \Leftrightarrow \mathbb{E}[X^2] = \text{Var}[X] + \mathbb{E}[X]^2$ ,

We have:

$$\mathbb{E}_{\mathcal{Z}_k} [\|g(w_k; \mathcal{Z}_k)\|^2] = \|\mathbb{E}_{\mathcal{Z}_k} [g(w_k; \mathcal{Z}_k)]\|_2^2 + \text{Var}_{\mathcal{Z}_k} [g(w_k; \mathcal{Z}_k)] \quad (2.2.i)$$

Since  $0 \leq \|\mathbb{E}_{\mathcal{Z}_k} [g(w_k; \mathcal{Z}_k)]\|_2 \leq C_0 \|\nabla f(w_k)\|_2$ , we have:

$$\|\mathbb{E}_{\mathcal{Z}_k} [g(w_k; \mathcal{Z}_k)]\|_2^2 \leq C_0^2 \|\nabla f(w_k)\|_2^2. \quad (2.2.ii)$$

Substituting (2.2.ii) and  $\text{Var}_{\mathcal{Z}_k} [g(w_k; \mathcal{Z}_k)] \leq M + M_v \|\nabla f(w_k)\|_2^2$  into (2.2.i) gives:

$$\begin{aligned} \mathbb{E}_{\mathcal{Z}_k} [\|g(w_k; \mathcal{Z}_k)\|^2] &\leq C_0^2 \|\nabla f(w_k)\|_2^2 + M + M_v \|\nabla f(w_k)\|_2^2 = \\ &= M + (C_0 + M_v) \|\nabla f(w_k)\|_2^2 \end{aligned}$$

We have thus proved that  $\mathbb{E}_{\mathcal{Z}_k} [\|g(w_k; \mathcal{Z}_k)\|^2] \leq \alpha + \beta \|\nabla f(w_k)\|_2^2$

and found that  $\alpha = M$ ,  $\beta = C_0 + M_v$ .

Q.E.D.



HW2.3 For SGD with non-convex objective functions. Prove that with square summable but not summable step-size, for any  $k \in \mathbb{N}$ .

$$\mathbb{E} \left[ \sum_{k \in [k]} \alpha_k \|\nabla f(w_k)\|_2^2 \right] < \infty.$$

and therefore,  $\mathbb{E} \left[ \frac{1}{\sum_{k \in [k]} \alpha_k} \sum_{k \in [k]} \alpha_k \|\nabla f(w_k)\|_2^2 \right] \xrightarrow{k \rightarrow \infty} 0$

Proof:

Generic SG on  $L$ -smooth function satisfies:

$$\mathbb{E}[f(w_{k+1})] - f(w_k) \leq -\left(c - \frac{1}{2}\alpha_k L M_G\right) \alpha_k \|\nabla f(w_k)\|_2^2 + \frac{1}{2} \alpha_k^2 L M$$

(1)

For  $k \in [k]$ , the total expectation of (1) is:

$$\begin{aligned} f_{\min} - f(w_1) &\leq \mathbb{E}[f(w_{k+1})] - f(w_1) \\ &\leq -c \mathbb{E} \left[ \sum_{k \in [k]} \alpha_k \|\nabla f(w_k)\|_2^2 \right] + \frac{1}{2} L M_G \mathbb{E} \left[ \sum_{k \in [k]} \alpha_k^2 \|\nabla f(w_k)\|_2^2 \right] \\ &\quad + \frac{L M}{2} \sum_{k \in [k]} \alpha_k^2 \end{aligned}$$

Then,

$$\mathbb{E} \left[ \sum_{k \in [k]} \alpha_k \|\nabla f(w_k)\|_2^2 \right] \leq \frac{f(w_1) - f_{\min}}{c} + \frac{L M_G}{2c} \mathbb{E} \left[ \sum_{k \in [k]} \alpha_k^2 \|\nabla f(w_k)\|_2^2 \right] + \frac{L M}{2c} \sum_{k \in [k]} \alpha_k^2$$

(2)

Since  $f$  is  $L$ -smooth,  $f$  must be Lipschitz continuous,

then  $\|\nabla f(w_k)\|$  is bounded on  $\{w_k: \|w_k\| \leq D\}$ , i.e.,  $\|\nabla f(w_k)\| \leq B$ .

Hence,

$$\mathbb{E} \left[ \sum_{k \in [k]} \alpha_k^2 \|\nabla f(w_k)\|_2^2 \right] \leq B^2 \mathbb{E} \left[ \sum_{k \in [k]} \alpha_k^2 \right] < \infty.$$

Since in (2),  $\frac{f(w_1) - f_{\min}}{c}$ ,  $\frac{L M_G}{2c} \mathbb{E} \left[ \sum_{k \in [k]} \alpha_k^2 \|\nabla f(w_k)\|_2^2 \right]$  (since  $\sum_{k \in [k]} \alpha_k^2 < \infty$ ),  $\frac{L M}{2c} \sum_{k \in [k]} \alpha_k^2$  are all bounded,

Then we have  $\mathbb{E} \left[ \sum_{k \in [k]} \alpha_k \|\nabla f(w_k)\|_2^2 \right] < \infty$ .



From (2) divide each by  $\sum_{k \in \mathcal{K}} \partial_k$  we have

$$\mathbb{E} \left[ \frac{1}{\sum \partial_k} \sum \partial_k \|\nabla f(w_k)\|^2 \right] \leq \frac{f(w_1) - f_{\inf}}{c \cdot \sum \partial_k} + \frac{2M_0 B^2 + 4M}{2c} \cdot \frac{\sum \partial_k^2}{\sum \partial_k}$$

since  $\lim_{k \rightarrow \infty} \sum_{k \in \mathcal{K}} \partial_k = \infty$  and  $\lim_{k \rightarrow \infty} \sum_{k \in \mathcal{K}} \partial_k^2 < \infty$ ,

it holds that,

$$\frac{1}{\sum \partial_k} \xrightarrow{k \rightarrow \infty} 0 \quad \text{and} \quad \frac{\sum \partial_k^2}{\sum \partial_k} \rightarrow 0$$

Therefore,

$$\mathbb{E} \left[ \frac{1}{\sum \partial_k} \sum \partial_k \|\nabla f(w_k)\|^2 \right] \xrightarrow{k \rightarrow \infty} 0$$